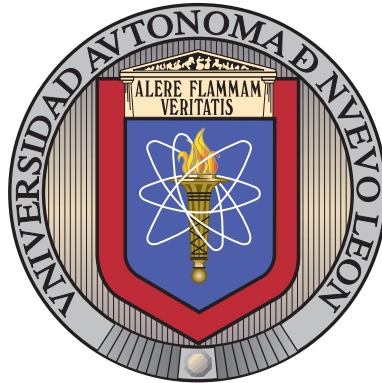


UNIVERSIDAD AUTÓNOMA DE NUEVO LEÓN

FACULTAD DE INGENIERÍA MECÁNICA Y ELÉCTRICA

SUBDIRECCIÓN ACADÉMICA



MODELADO Y VISUALIZACIÓN DE RELACIONES
ENTRE CONTAMINANTES DEL AIRE Y SALUD
PÚBLICA

POR

SELENE BERENICE PRADO PRADO

COMO REQUISITO PARCIAL PARA OBTENER EL GRADO DE
INGENIERÍA EN TECNOLOGÍA DE SOFTWARE

JULIO 2022

UNIVERSIDAD AUTÓNOMA DE NUEVO LEÓN

FACULTAD DE INGENIERÍA MECÁNICA Y ELÉCTRICA

SUBDIRECCIÓN ACADÉMICA



MODELADO Y VISUALIZACIÓN DE RELACIONES
ENTRE CONTAMINANTES DEL AIRE Y SALUD
PÚBLICA

POR

SELENE BERENICE PRADO PRADO

COMO REQUISITO PARCIAL PARA OBTENER EL GRADO DE
INGENIERÍA EN TECNOLOGÍA DE SOFTWARE

JULIO 2022

Universidad Autónoma de Nuevo León
Facultad de Ingeniería Mecánica y Eléctrica
Subdirección Académica

Los miembros del Comité de Tesis recomendamos que la Tesis «Modelado y visualización de relaciones entre contaminantes del aire y salud pública», realizada por el alumno Selene Berenice Prado Prado, con número de matrícula 1810042, sea aceptada para su defensa como requisito parcial para obtener el grado de Ingeniería en Tecnología de Software.

El Comité de Tesis

Dra. Satu Elisa Schaeffer

Coasesora

Dra. Sara Elena Garza Villarreal

Coasesora

Dr. José Arturo Berrones Santos

Revisor

Dr. Romeo Sánchez Nigenda

Revisor

Vo. Bo.

Dr. Fernando Banda Muñoz

Subdirección Académica

San Nicolás de los Garza, Nuevo León, julio 2022

ÍNDICE GENERAL

Agradecimientos	IX
Resumen	X
1. Introducción	1
1.1. Motivación	3
1.2. Hipótesis	3
1.3. Objetivos	3
1.3.1. Objetivo general	3
1.3.2. Objetivos específicos	4
2. Antecedentes	5
2.1. Monitoreo de calidad del aire	5
2.2. Series de tiempo	7
2.3. Clasificación de enfermedades	8
2.4. Regresión lineal	8
2.5. Regresión lineal múltiple	8

3. Estado del arte	9
3.1. Trabajos relacionados	10
3.2. Comparación de trabajos	12
3.2.1. Áreas de oportunidad	14
4. Solución propuesta	15
4.1. Diseño de la solución propuesta	15
4.1.1. Recolección de datos	16
4.1.2. Selección y agrupación de datos	16
4.1.3. Visualización de la evolución de las variables	17
4.1.4. Implementación de modelos	19
4.2. Implementación de la solución propuesta	20
5. Experimentos	27
5.1. Diseño experimental	27
5.1.1. Datos de entrada	27
5.1.2. Visualización de datos	28
5.1.3. Generación de modelos	29
5.2. Resultados	31
5.2.1. Experimento A: Datos de niveles de PM10	31
5.2.2. Experimento B: Niveles de PM2.5	38
5.3. Discusión	44

ÍNDICE GENERAL	VI
6. Conclusiones	46
6.1. Contribuciones	47
6.2. Trabajo a futuro	48
A. CIE y sus nombres de enfermedades	49

ÍNDICE DE FIGURAS

1.1. Localización de las estaciones de monitoreo de la calidad del aire. . .	2
4.1. Ejemplo de gráfico de radar	18
4.2. Fases del desarrollo de la solución.	21
5.1. Series de tiempo 2017 PM10 y O809	32
5.2. Correlaciones 2017 PM10	33
5.3. Series de tiempo 2018 PM10 y O809	35
5.4. Correlaciones 2018 PM10	37
5.5. Series de tiempo 2017 PM2.5 y O809	38
5.6. Correlaciones 2017 PM2.5	40
5.7. Series de tiempo 2018 PM2.5 y O809	41
5.8. Correlaciones 2018 PM2.5	43

ÍNDICE DE CUADROS

1.1. Coordenadas de las estaciones de monitoreo de la calidad del aire. . .	2
3.1. Comparación de trabajos	13
4.1. Herramientas utilizadas.	16
5.1. Resultados obtenidos PM10 2017	32
5.2. Resultados regresión lineal múltiple PM10 2017	34
5.3. Resultados obtenidos PM10 2018	36
5.4. Resultados regresión lineal múltiple PM10 2018	36
5.5. Resultados obtenidos PM2.5 2017	39
5.6. Resultados regresión lineal múltiple PM2.5 2017	39
5.7. Resultados obtenidos PM2.5 2018	42
5.8. Resultados regresión lineal múltiple PM2.5 2018	42
5.9. Especificaciones técnicas del equipo de cómputo	45
A.1. CIE mencionadas en los Experimentos y el nombre de la enfermedad.	49

AGRADECIMIENTOS

Quiero agradecer a la Dra. Elisa por el apoyo durante el desarrollo de mi tesis y por la motivación y conocimientos brindados para seguir desarrollándome profesionalmente en lo que me gusta. Al programa PAICYT-UANL por su contribución brindada bajo las claves CE1421-20 y CE1842-21.

A mis padres, Lilia Prado López y Adan Alfaro Lerma, por su apoyo y motivación constante desde siempre. A mis hermanos Angel, Estrella, y Adali, a quienes he visto crecer y de quienes he aprendido mucho.

RESUMEN

Selene Berenice Prado Prado.

Candidato para obtener el grado de Ingeniería en Tecnología de Software.

Universidad Autónoma de Nuevo León.

Facultad de Ingeniería Mecánica y Eléctrica.

Título del estudio: MODELADO Y VISUALIZACIÓN DE RELACIONES ENTRE CONTAMINANTES DEL AIRE Y SALUD PÚBLICA.

Número de páginas: 52.

OBJETIVOS Y MÉTODO DE ESTUDIO: El objetivo de la investigación es generar modelos que permitan visualizar relaciones entre contaminantes atmosféricos y salud pública. Los modelos generados se utilizan en conjunto con datos obtenidos de la Secretaría de Salud del Gobierno de México y registros de los niveles de los contaminantes presentes en el área metropolitana de Monterrey.

El tener un modelo que permita visualizar relaciones entre contaminantes atmosféricos y salud pública que sea utilizado con datos confiables y verídicos pueden ayudar a entender el impacto que tiene el aumento del nivel de contaminantes atmosféricos.

CONTRIBUCIONES Y CONCLUSIONES: En la investigación se exploran diversas maneras de visualizar la información, además de generar modelos que permiten analizar las relaciones existentes entre los niveles de determinados contaminantes y CIE. Las visualizaciones generadas son series de tiempo y gráficos de radar, y modelos de regresión lineal y de regresión lineal múltiple.

Firma de la coasesora: _____

Dra. Satu Elisa Schaeffer

Firma de la coasesora: _____

Dra. Sara Elena Garza Villarreal

INTRODUCCIÓN

El *aprendizaje máquina*¹ es un área dentro de la *ciencia de datos*² que puede ayudar a crear dichos modelos para tener una más eficiente visualización cuando se trabaja con una gran cantidad de datos, que es lo que se requiere para el presente trabajo. El área de la ciencia de datos es muy útil ya que permite trabajar con grandes cantidades de datos aminorando la cantidad de tiempo empleado en la creación de gráficos que permitan visualizar los datos. El crear modelos para la visualización de datos ayuda a observar con mayor claridad los datos para encontrar relaciones entre ellos.

La tarea en el presente proyecto es utilizar modelos para visualizar las relaciones entre los contaminantes del aire y salud pública. Para la realización de los experimentos se tienen datos de ingresos hospitalarios provenientes de la base de datos de la Secretaría de Salud del Gobierno de México [9]. También se tienen registros de los niveles de algunos contaminantes del aire presentes en el área metropolitana de Monterrey, capturados por las estaciones de monitoreo pertenecientes al Sistema Integral de Monitoreo Ambiental (SIMA) [20] mostradas en la figura 1.1 y en el cuadro 1.1.

¹Traducido como *machine learning* en inglés, tiene como objetivo desarrollar técnicas que les permitan a las computadoras aprender.

²Traducido como *data science* en inglés, involucra métodos para extraer conocimiento de datos, eso con la finalidad de que haya un mejor entendimiento de los datos.

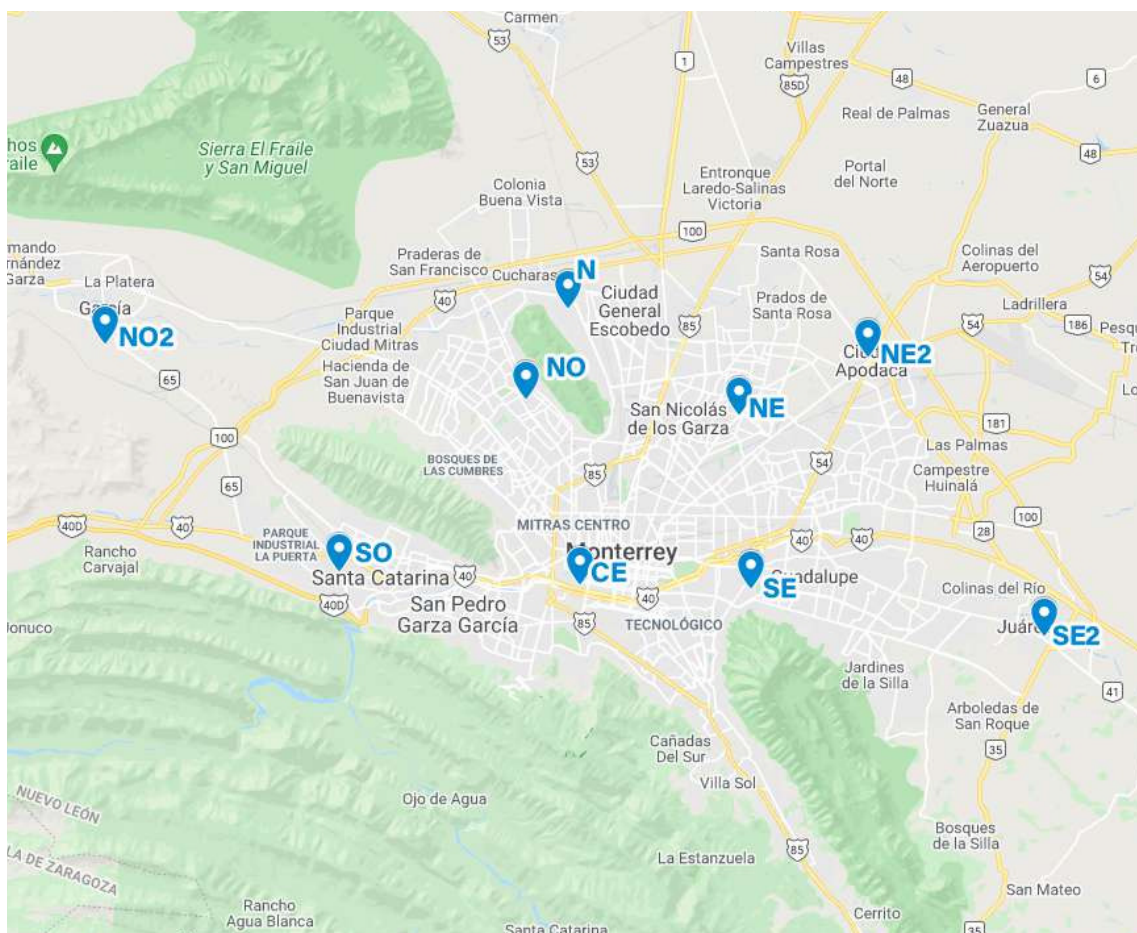


Figura 1.1: Localización de las estaciones de monitoreo de la calidad del aire.

Cuadro 1.1: Coordenadas de las estaciones de monitoreo de la calidad del aire.

Estación	Latitud	Longitud
SE	25.668	-100.249
NE	25.750	-100.255
CE	25.670	-100.338
NO	25.757	-100.366
SO	25.676	-100.464
NO2	25.783	-100.586
NE2	25.777	-100.188
N	25.800	-100.344
SE2	25.646	-100.096

1.1 MOTIVACIÓN

Existen investigaciones que ya han estudiado las relaciones entre contaminantes del aire y salud pública, sin embargo, con el presente trabajo se busca aportar a la creación de nuevas herramientas que permitan observar y estudiar dichas relaciones con el fin de ayudar a tomar medidas adecuadas que permitan aminorar los efectos negativos de los contaminantes del aire en la salud.

1.2 HIPÓTESIS

Los modelos de regresión permiten obtener gráficos donde se pueden observar las relaciones entre el número de ingresos hospitalarios y los niveles de contaminantes del aire.

1.3 OBJETIVOS

En esta sección se establece el objetivo general y los objetivos específicos sobre los que se orienta el presente trabajo.

1.3.1 OBJETIVO GENERAL

Generar, implementar y evaluar modelos que muestran las relaciones existentes entre contaminantes del aire y salud pública tiene la finalidad de apoyar a la implementación de estrategias que aminoran los efectos negativos de los contaminantes del aire en la salud de las personas. Con los modelos generados se puede tener una herramienta que permite identificar gráficamente las relaciones con solo proporcionarle el conjunto de datos.

39 1.3.2 OBJETIVOS ESPECÍFICOS

- 40 ■ Generar, implementar y evaluar modelos de regresión que permiten cuantifi-
41 car las relaciones entre contaminantes del aire y salud pública a partir de un
42 conjunto de datos.

- 43 ■ Diseñar e implementar visualizaciones interactivas que permiten explorar los
44 modelos implementados y su validez estadística.

- 45 ■ Evaluar la eficacia de los modelos generados para tener noción de la fiabilidad
46 del análisis realizado a partir de los resultados que tales modelos producen.

ANTECEDENTES

49 Existen factores ambientales que afectan la salud de una comunidad como:
50 el abastecimiento de agua potable y el saneamiento, la vivienda y el hábitat, la
51 alimentación, la contaminación ambiental, el empleo de productos químicos y los
52 riesgos ocupacionales [6].

53 Contaminación del aire es un término usado para describir la presencia de uno
54 o más contaminantes en la atmósfera, cuyas cantidades y características pueden re-
55 sultar perjudiciales o interferir con la salud, el bienestar u otros procesos ambientales
56 naturales [7].

57 En este capítulo se presentan los fundamentos y definiciones de los conceptos
58 más relevantes para el tema de estudio abordado.

2.1 MONITOREO DE CALIDAD DEL AIRE

60 Existen diversos estudios que muestran la existencia de potenciales efectos a la
61 salud cuando en el aire están presentes contaminantes en forma de partículas, gases
62 o agentes biológicos.

63 Korc y Sáenz [15] mencionan que desde inicios de 1950 se observa una preocu-
64 pación por los contaminantes del aire en los países de América Latina y el Caribe.
65 Las universidades y dependencias de los ministerios de salud fueron los organismos
66 que realizaron las primeras mediciones de contaminación en el aire.

67 En 1965, el Consejo Directivo de la Organización Panamericana de la Salud
68 (OPS) recomendó el establecer programas de investigación de la contaminación del
69 agua y del aire, con el objetivo de colaborar en el desarrollo de políticas adecuadas
70 de control [12].

71 Mediante el Centro Panamericano de Ingeniería Sanitaria y Ciencias del Am-
72 biente (CEPIS), la OPS acordó establecer una red de estaciones de muestreo de la
73 contaminación del aire. En junio de 1967 la Red Panamericana de Muestreo Nor-
74 malizado de la Contaminación del Aire (REDPANAIRE) inició sus operaciones re-
75 colectando muestras diarias de partículas totales en suspensión (PTS) y de SO_2 . La
76 REDPANAIRE comenzó con ocho estaciones y a fines de 1973 tenía un total de 88
77 estaciones distribuidas en 26 ciudades de 14 países [12].

78 Para diciembre de 1973 se habían recolectado más de 350,000 datos sobre la
79 calidad del aire, en los que se observa que algunas ciudades mostraban una tendencia
80 al incremento de los niveles de contaminación [12].

81 En 1980 la REDPANAIRE desapareció y pasó a formar parte del Programa
82 Global de Monitoreo de la Calidad del Aire, iniciado en 1976 por la OMS y el
83 Programa de las Naciones Unidas para el Medio Ambiente (PNUMA), como parte
84 de un sistema global de monitoreo ambiental llamado GEMS por sus siglas en inglés
85 *Global Environmental Monitoring System*.

86 En la década de 1990, la OMS instituyó, con carácter global, el Sistema de
87 Información para el Control de la Calidad del Aire llamado AMIS por sus siglas en
88 inglés *Air Management Information System*. Entre las actividades más destacadas
89 de AMIS se incluye el coordinar las bases de datos sobre temas relacionados con la
90 calidad del aire.

91 En Nuevo León, México, las operaciones de la Red Automática de Monitoreo
92 Atmosférico iniciaron en 1993. Dicha red en sus inicios contaba con cinco estaciones
93 fijas de monitoreo continuo de monóxido de carbono (CO), dióxido de azufre (SO₂),
94 óxidos de nitrógeno (NO_x), ozono (O₃) y partículas de tamaño menor a 10 micróme-
95 tros (PM10) [15]. Como se muestra en la figura 1.1 y en el cuadro 1.1, actualmente
96 se cuenta con nueve estaciones fijas.

97 2.2 SERIES DE TIEMPO

98 Korc y Sáenz [15] mencionan que las relaciones entre niveles de concentraciones
99 de contaminantes del aire y los efectos sobre la salud generalmente son obtenidas
100 de estudios epidemiológicos de series de tiempo. Uno de los diseños epidemiológicos
101 más utilizados son los estudios de series temporales. Con esos diseños se analizan
102 las variaciones en el tiempo de la exposición al contaminante y el indicador de salud
103 estudiado en una población [2].

104 Las series de tiempo se pueden definir como un conjunto de observaciones *ot*
105 tomadas en un tiempo *t* determinado. Los estudios de series de tiempo relacionan
106 estadísticamente los cambios temporales en la repercusión de cambios en la concen-
107 tración de un contaminante en la población [5].

108 Para mostrar datos en una serie de tiempo, especialmente en el área médica,
109 estos suelen agruparse en *semanas epidemiológicas*¹. El agrupamiento en semanas
110 epidemiológicas a diferencia del agrupamiento traducido como *clustering* en inglés
111 que consiste en crear grupos de objetos con similitudes entre ellos a partir de datos
112 no etiquetados, consiste en agrupar los objetos en semanas del año dependiendo de
113 la fecha en que se registró el objeto.

¹Una semana epidemiológica es un estándar de medición temporal que se utiliza para comparar datos en ventanas de tiempo definidas. La primera semana epidemiológica del año termina el primer sábado de enero de cada año [1].

2.3 CLASIFICACIÓN DE ENFERMEDADES

Existe un instrumento estadístico y sanitario para identificar enfermedades llamado Clasificación Internacional de Enfermedades (CIE), cuya finalidad es entender las causas de morbilidad y mortalidad de la población y así mejorar la calidad de vida de la misma. Es en base a un criterio epidemiológico y sanitario establecido por Farr a finales del siglo XIX que esta clasificación agrupa enfermedades en epidémicas, generales, locales ordenadas por origen geográfico, trastornos del desarrollo y lesiones [18]. Para lograr distinguirlas se emplea un código alfanumérico que consiste de una letra en la primera posición, seguida de dos dígitos, un punto decimal y un último dígito. El rango de valores va de A00.0 a Z99.9.

2.4 REGRESIÓN LINEAL

La tendencia w_0 de una serie de tiempo puede ser obtenida a partir de una regresión lineal de la misma [8]. Una regresión lineal es una metodología inferencial supervisada que busca predecir valores y dado un vector de variables de entrada t por medio del ajuste de coeficientes w de la función lineal

$$\hat{y}(t, w) = w_0 + w_1x_1 + \dots + w_tx_t. \quad (2.1)$$

2.5 REGRESIÓN LINEAL MÚLTIPLE

Un modelo de regresión múltiple es un modelo complemento de la regresión lineal simple, el cual tiene dos o más variables independientes k que pueden influir en una variable dependiente y . Peniche-Camps y Cortez-Huerta [19] expresan la regresión múltiple mediante la siguiente ecuación:

$$y = \beta_0 + \beta_1x_1 + \dots + \beta_kx_k + \varepsilon. \quad (2.2)$$

ESTADO DEL ARTE

136 En el presente capítulo se estudia literatura reciente relacionada con el presente
137 trabajo, esto con el objetivo de revisar distintos métodos para resolver el problema
138 planteado en el presente trabajo y, además, también revisar implementaciones simi-
139 lares para resolver problemas distintos. Lo anterior tiene la finalidad de comparar
140 los trabajos revisados e identificar áreas de oportunidad en ellos.

141 En la primera sección, *trabajos relacionados*, se recopilan obras con caracterís-
142 ticas relacionadas al presente trabajo, ya sean relacionados con el problema que se
143 pretende resolver o con los métodos empleados para buscar su resolución.

144 En la segunda sección, *análisis comparativo*, se comparan las distintas caracte-
145 rísticas de los trabajos revisados, de esa forma se pueden determinar las principales
146 ventajas y desventajas de cada trabajo.

147 Finalmente, en la tercera sección, *áreas de oportunidad*, se realiza una conclu-
148 sión acerca de los resultados obtenidos del análisis comparativo.

3.1 TRABAJOS RELACIONADOS

Se recopila literatura relacionada desde el año 2017 hasta el año 2021. En esta sección los trabajos se mencionan en orden cronológico tomando en cuenta su año de publicación.

Martín y Bayle [17] estudian la relación entre los niveles de contaminantes ambientales y la presencia de casos de enfermedades respiratorias en las consultas pediátricas. La variable dependiente analizada es la demanda en las consultas pediátricas por bronquiolitis, episodios de broncoespasmo y procesos respiratorios de vías altas. Como variables independientes se tienen los valores de contaminación ambiental. Se calculan coeficientes de correlación y regresión lineal múltiple.

Guarnaccia *et al.* [10] abordan la necesidad de monitoreo, control y predicción de la pendiente de los niveles de contaminantes del aire. Para abordar el problema de investigación utilizan modelos ARIMA.

Julia *et al.* [13] estudian la asociación entre la exposición a largo plazo a la contaminación del aire y la metilación del ADN. Para ello realizan un estudio utilizando modelos de regresión lineal robustos para analizar la asociación entre la exposición al NO₂ y a las partículas PM10 y PM2.5.

Zhang *et al.* [22] en su estudio abordan los niveles de contaminación del aire y su asociación con la presencia de presión sanguínea elevada en niños y adolescentes. La exposición a partículas PM10 y PM2.5 son estimadas con un modelo espacio-temporal. Son utilizados modelos lineales de efectos mixtos y modelos de regresión logística para investigar la asociación entre la exposición a partículas PM y presión sanguínea e hipertensión.

Kim *et al.* [14] estudian la relación entre los niveles de contaminación del aire y la obesidad y problemas cardiometabólicos. Para dicho estudio emplean modelos de regresión lineal.

175 Liu *et al.* [16] examinan las asociaciones entre la exposición temprana a la
176 contaminación del aire y la incidencia de asma y rinitis alérgica desde el nacimiento
177 hasta la adolescencia. Para su estudio utilizan modelos de regresión.

178 To *et al.* [21] estudian la asociación entre la exposición temprana a los contami-
179 nantes del aire y los egresos hospitalarios por asma. Para su estudio aplican modelos
180 de regresión logística para el análisis de datos.

181 Breton *et al.* [4] abordan el estudio de la relación entre los niveles de contami-
182 nación del aire y el número de admisiones hospitalarias. Para ello se construye un
183 modelo basado en la distribución de Poisson.

184 Gupta *et al.* [11] estudian la relación entre la mortalidad del coronavirus
185 (COVID-19) y la contaminación del aire. Para dicho estudio emplean un modelo de
186 regresión lineal para establecer la relación entre los parámetros de la contaminación
187 del aire (concentraciones de PM10 o PM2.5) y la variable de respuesta (porcentaje
188 de mortalidad por unidad de casos reportados).

3.2 COMPARACIÓN DE TRABAJOS

189

190 La mayoría de los trabajos encontrados emplean modelos de regresión lineal
191 o modelos de predicción. Además, en todos los trabajos encontrados el problema
192 tratado presenta una alta relación con el problema abordado en el presente trabajo
193 de tesis. El análisis comparativo de los trabajos relacionados se hace en base de los
194 siguientes puntos:

195 **Modelos de regresión lineal:** Son aquellos que ayudan a estudiar la relación en-
196 tre una variable dependiente y una o más variables independientes.

197 **Modelos de predicción:** Son aquellos que ayudan a hacer predicciones de una
198 variable.

199 **Evaluación de modelos:** Se refiere a la utilización de técnicas para evaluar la
200 eficacia de los modelos generados.

201 **Estudio de contaminantes del aire:** Se refiere a que el tema de estudio incluya
202 uno o más contaminantes del aire.

203 **Estudio de problemas de salud:** Se refiere a que el tema de estudio incluya uno
204 o más problemas de salud.

205 En el cuadro 3.1 se desglosan las características presentes que se pueden en-
206 contrar en las investigaciones citadas y su relación con la investigación con la que se
207 está trabajando actualmente.

Cuadro 3.1: Comparación de trabajos frente al desarrollado, donde ✓ indica que cumple con esta característica y × no cumple con esta característica.

Trabajo	Modelos de regresión lineal	Modelos de predicción	Evaluación de modelos	Estudio de contaminantes del aire	Estudio de problemas de salud
Martín y Bayle [17]	✓	×	×	✓	✓
Guarnaccia <i>et al.</i> [10]	×	✓	✓	✓	×
Julia <i>et al.</i> [13]	✓	✓	×	✓	✓
Zhang <i>et al.</i> [22]	✓	✓	×	✓	✓
Kim <i>et al.</i> [14]	✓	×	×	✓	✓
Liu <i>et al.</i> [16]	×	✓	✓	✓	✓
To <i>et al.</i> [21]	×	×	×	✓	✓
Breton <i>et al.</i> [4]	✓	✓	×	✓	✓
Gupta <i>et al.</i> [11]	✓	✓	×	✓	✓
El presente trabajo	✓	✓	✓	✓	✓

208 3.2.1 ÁREAS DE OPORTUNIDAD

209 Como se puede observar en el cuadro 3.1, la mayoría de los trabajos encontrados
210 abordan el estudio de los contaminantes del aire y salud con excepción de Guarnaccia
211 *et al.* [10] que se enfocan en la predicción de niveles de contaminantes del aire, lo
212 cual puede indicar que la relación entre los contaminantes del aire y salud es un tema
213 de relevancia en la actualidad.

214 Ya que la mayoría de los trabajos encontrados estudian la relación entre conta-
215 minantes del aire y salud, la mayoría de los trabajos emplean modelos de regresión
216 lineal por que es una buena opción para el estudio de relaciones entre variables. Las
217 excepciones, además de la ya anteriormente mencionada, son Liu *et al.* [16] y To *et*
218 *al.* [21] quienes emplean otros tipos de modelos de regresión.

219 En el presente trabajo se elaboran modelos de predicción para el tratamiento
220 de los datos empleados para los experimentos, ya que como mencionan Zhang *et al.*
221 [22], una de las limitaciones en este tipo de estudios es los campos sin llenar en los
222 registros de datos.

223 En el presente trabajo también se emplean técnicas para evaluar los modelos
224 generados. Solo en tres de los trabajos encontrados se aborda la evaluación de los
225 modelos empleados, y al ser incluida en el presente estudio, puede representar una
226 distinción.

227

CAPÍTULO 4

228

SOLUCIÓN PROPUESTA

229

En este capítulo se comparte la propuesta de diseño de la solución para el
230 problema de investigación abordado en este presente trabajo, así como su imple-
231 mentación.

232

4.1 DISEÑO DE LA SOLUCIÓN PROPUESTA

233

En diseño de la solución propuesta se plantean las herramientas utilizadas y
234 los pasos seguidos para la solución propuestas.

235

Las herramientas utilizadas en la presente investigación se muestran en el cua-
236 dro 4.1.

Cuadro 4.1: Herramientas utilizadas.

Herramienta	Versión	URL
Python	3.8.8	https://www.python.org/
Jupyter Notebook	6.3.0	https://jupyter.org/
Imageio	2.9.0	https://imageio.readthedocs.io/
Latextable	0.2.1	https://pypi.org/project/latextable/
Matplotlib	3.3.4	https://matplotlib.org/
NumPy	1.20.1	http://www.numpy.org/
Pandas	1.2.4	https://pandas.pydata.org/
Seaborn	0.11.1	https://seaborn.pydata.org/
Scikit-learn	0.24.1	https://scikit-learn.org/
SciPy	1.6.2	https://docs.scipy.org/
Statsmodels	0.12.2	https://www.statsmodels.org/
Texttable	1.6.4	https://pypi.org/project/texttable/

237 4.1.1 RECOLECCIÓN DE DATOS

238 La primera fase es la recolección de datos. El objetivo es tener un archivo que
 239 contenga datos de los niveles de uno o más contaminantes del aire en años recientes
 240 y también del mismo lugar contar con datos del número de egresos hospitalarios
 241 durante esos años.

242 4.1.2 SELECCIÓN Y AGRUPACIÓN DE DATOS

243 Después de la recolección de datos se seleccionan cuáles van a ser utilizados
 244 para los experimentos. Para ello se utiliza **Python** con la librería **Pandas** que per-
 245 mite la manipulación de datos. Para la selección y agrupación de datos en semanas
 246 epidemiológicas se sigue el procedimiento mostrado en el algoritmo 1.

Algoritmo 1 Selección y agrupamiento de datos

```

1:  $a \leftarrow$  años de los que se obtuvieron datos
2: for  $i \in a$  do
3:    $contaminantes \leftarrow$  nombre del archivo .csv que contiene los datos de los
     contaminantes en el año  $i$ 
4:   Leer en  $contaminantes$  las columnas fecha y contaminante
5:    $egresos \leftarrow$  nombre del archivo .csv que contiene los datos de los contaminan-
     tes en el año  $i$ 
6:   Leer en  $contaminantes$  las columnas fecha, padecimiento y estado
7:    $estado \leftarrow$  estado del que se quieren obtener datos
8:   Seleccionar en  $contaminantes$  los datos del  $estado$ 
9: end for

```

247 4.1.3 VISUALIZACIÓN DE LA EVOLUCIÓN DE LAS VARIABLES

248 Al ya tener seleccionados los datos a utilizar se procede a elaborar gráficos en
249 Python que muestran la evolución de las variables en el tiempo. Para ello se generan
250 los tipos de gráficos discutidos a continuación.

251 4.1.3.1 SERIES DE TIEMPO

252 Se realizan series de tiempo en Python con ayuda de la librería Matplotlib,
253 Scikit-learn y Seaborn, ya que son herramientas accesibles que ayudan a la gene-
254 ración de este tipo de gráficos.

255 4.1.3.2 GRÁFICOS DE RADAR

256 Los gráficos de radar o diagramas de telaraña son otra manera de visualizar
257 un conjunto de datos. Sirven para representar datos en un diagrama bidimensional
258 cuando hay múltiples variables y para comparar variables visualizando si existen
259 valores o patrones de evolución en el tiempo similares entre ellas. Se puede utilizar
260 para graficar un grupo de variables o para comparar varios grupos de las mismas
261 variables. Cada variable se grafica en ejes que se organizan equitativamente alrededor

de un punto de origen central. Las líneas de cuadrícula conectan los ejes y al conectar cada grupo de variables en el gráfico, las conexiones forman un polígono.

Es por ello que en el presente trabajo se elaboran gráficos de radar con ayuda de Python y las librerías NumPy y Matplotlib. En la figura 4.1 se muestra un ejemplo de los gráficos de telaraña generados. Para su interpretación: se identifica la categoría que representa cada eje, se determina cómo se relaciona cada categoría con las demás, se observa la forma completa creada, se lee alrededor de la rueda y se comparan los datos.

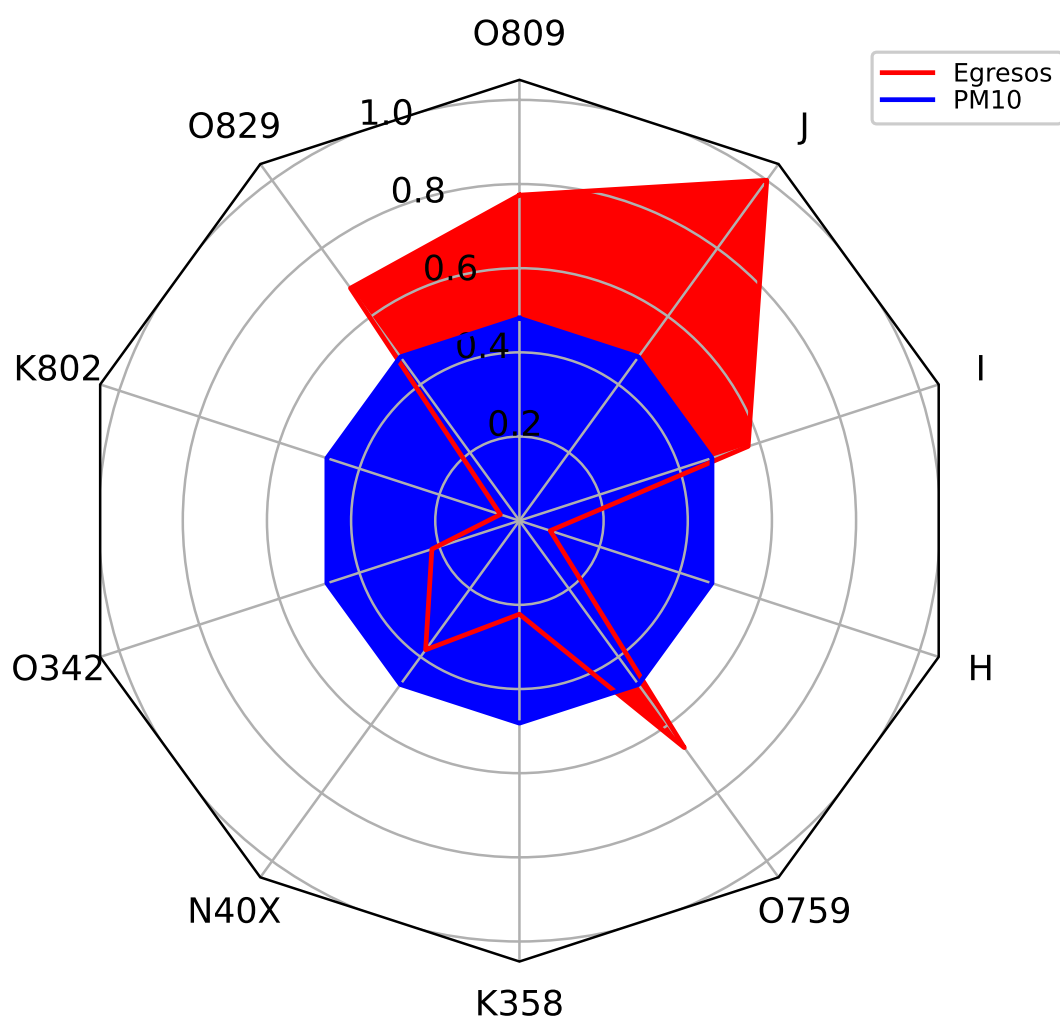


Figura 4.1: Nivel del contaminante PM10 y egresos por CIE en la semana 1 del año 2018.

270 4.1.4 IMPLEMENTACIÓN DE MODELOS

271 Después de haber generado gráficos para la visualización de la evolución de
272 las variables, se procede a generar modelos para el estudio de la relación entre las
273 variables. Para ello se utiliza `Python` y la librería `Statsmodels`. Los tipos de modelos
274 generados se discuten a continuación.

275 4.1.4.1 REGRESIÓN LINEAL

276 Primeramente se calcula el coeficiente de correlación de Pearson y se verifica
277 que esté entre -1 y 1. Si un valor es cercano a 0, quiere decir que no hay dependencia
278 lineal. Si no hay una dependencia lineal no existe sustento para el modelo de regresión
279 lineal. El modelo de regresión lineal arroja un valor de R^2 que indica en qué grado
280 la variable independiente explica la varianza de la variable dependiente. Además,
281 se obtiene el valor p que indica la relevancia del resultado y se obtiene la raíz de
282 error cuadrático medio (RMSE) que indica cuántas unidades se alejan los valores
283 predichos por el modelo de los valores reales, eso ayuda a determinar el error del
284 modelo.

285 4.1.4.2 REGRESIÓN LINEAL MÚLTIPLE

286 En los modelos de regresión lineal múltiple se tiene más de una variable inde-
287 pendiente. Primero se calcula el coeficiente de correlación de Pearson y se verifica
288 que hay una dependencia lineal para generar el modelo de regresión lineal múltiple.
289 El modelo arroja un valor de R^2 que indica en qué grado las variables independien-
290 tes explican la varianza de la variable dependiente. Además, se obtiene el valor p
291 que indica la relevancia del resultado y se obtiene la raíz de error cuadrático medio
292 (RMSE) que indica cuántas unidades se alejan los valores predichos por el modelo
293 de los valores reales, eso ayuda a determinar el error del modelo.

294 4.2 IMPLEMENTACIÓN DE LA SOLUCIÓN PROPUESTA

295 En implementación de la solución propuesta se muestra el desarrollo reali-
296 zado de los puntos planteados en la sección 4.1. La figura 4.2 muestra las fases
297 seguidas para el desarrollo de la solución propuesta, la cual consiste en la genera-
298 ción de visualizaciones y modelos. El desarrollo del presente proyecto se encuen-
299 tra en el siguiente repositorio de Github: [https://github.com/selenebpradop/](https://github.com/selenebpradop/relaciones-contaminantes-salud/)
300 [relaciones-contaminantes-salud/](https://github.com/selenebpradop/relaciones-contaminantes-salud/).

301 En el fragmento de código 4.1 se muestra el proceso realizado para el procesa-
302 miento y agrupamiento de los datos en semanas epidemiológicas, esto para que los
303 datos puedan ser utilizados para generar figuras y modelos de regresión lineal.

304 El fragmento de código 4.2 muestra cómo es que se generan las series de tiempo,
305 esto después de haber procesado y agrupado los datos.

306 El fragmento de código mostrado en 4.3 genera una animación de gráficos de
307 radar al ingresarle como parámetros los datos ya procesados y agrupados. Dicha
308 animación es generada en formato .mp4 o .gif para poder visualizar la evolución de
309 las variables por semana del año y contaminante.

310 En el fragmento de código 4.4 se muestra cómo son generados los modelos de
311 regresión lineal después de obtener un coeficiente de correlación de Pearson entre -1
312 y 1.

313 El fragmento de código 4.5 muestra cómo se generan los modelos de regresión
314 lineal múltiple después de generar los modelos de regresión lineal individuales. Todos
315 los modelos generados por librería **Statsmodels** se guardan en formato .tex.

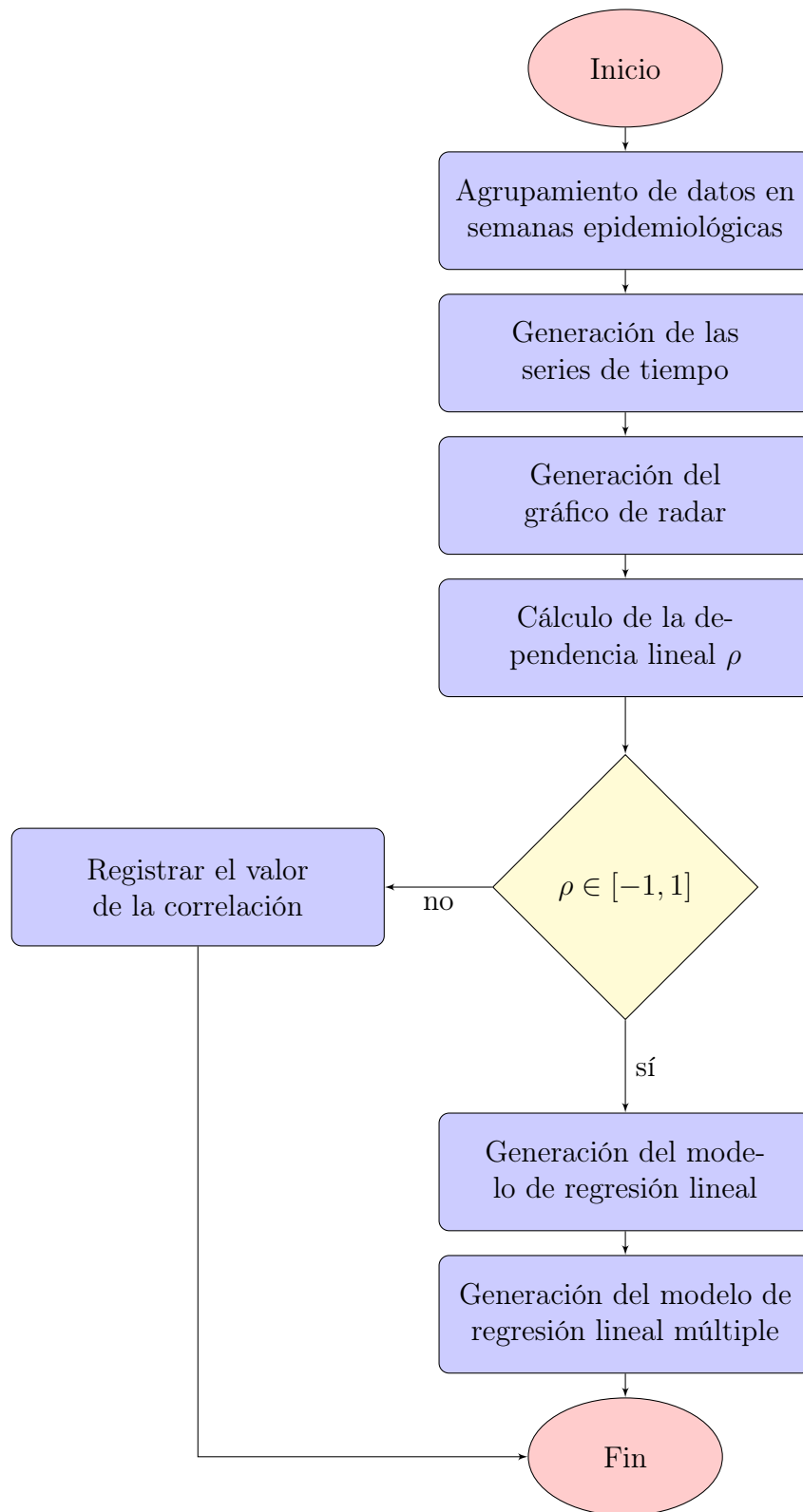


Figura 4.2: Fases del desarrollo de la solución.


```

1 316 import pandas as pd
2 317 from epiweeks import Week, date
3 318 from sklearn import preprocessing
4 319 import seaborn as sns
5 320 import matplotlib.pyplot as plt
6 321 import string
7 322
8 323 columns = ['timestamp', contaminante]
9 324 dataframec = pd.read_csv('filled.csv', usecols=columns).dropna()
10 325 strfddt = '%d-%b-%y %H'
11 326 dataframec['timestamp'] = pd.to_datetime(dataframec['timestamp'], errors = 'coerce', format=
12 327 strfddt)
13 328 dataframec = dataframec.dropna()
14 329 dataframec = dataframec.reset_index(drop=True)
15 330 dataframec['timestamp'] = dataframec['timestamp'].apply(lambda x: x.strftime('%Y-%m-%d %H'))
16 331 dataframeca = dataframec.loc[dataframec['timestamp'].str.startswith(año)]
17 332 dataframeca = dataframeca.reset_index(drop=True)
18 333 strfddt = '%Y-%m-%d %H'
19 334 dataframeca['timestamp'] = pd.to_datetime(dataframeca['timestamp'], errors = 'coerce', format=
20 335 strfddt)
21 336 dataframeca['sem'] = dataframeca['timestamp'].apply(lambda x: date(x.year, x.month, x.day))
22 337 dataframeca['sem'] = dataframeca['sem'].apply(lambda x: Week.fromdate(x))
23 338 dataframeca['sem'] = dataframeca['sem'].apply(lambda x: x.week)
24 339 columns = ['EGRESO', 'DIAG_INI']
25 340 csvegresos = 'EGRESO_' + año + '.csv'
26 341 dataframeea = pd.read_csv(csvegresos, usecols=columns).dropna()
27 342 dataframeea['EGRESO'] = pd.to_datetime(dataframeea['EGRESO'], errors = 'coerce', format=strfddt)
28 343 dataframeea = dataframeea.loc[dataframeea['ENTIDAD'] == entidad]
29 344 dataframeea = dataframeea.dropna()
30 345 dataframeea = dataframeea.reset_index(drop=True)
31 346 numañ = int(año)
32 347 dataframeea['sem'] = dataframeea['EGRESO'].apply(lambda x: date(x.year, x.month, x.day))
33 348 dataframeea['sem'] = dataframeea['sem'].apply(lambda x: Week.fromdate(x))
34 349 dataframeea['sem'] = dataframeea['sem'].apply(lambda x: x.week)
35 350 dataframeea['EGRESO'] = dataframeea['EGRESO'].apply(lambda x: x if (x.year==numañ) else pd.NaT)
36 351 dataframeea = dataframeea.dropna()
37 352 dataframeea = dataframeea.reset_index(drop=True)
38 353 dataframesca = pd.DataFrame()
39 354 dataframesca['sem'] = semanas.index
40 355 dataframesca[contaminante] = ''
41 356 n = len(semanas.index)
42 357 for i in range(n):
43 358     registrossem = dataframeca.loc[dataframeca['sem'] == i+1]
44 359     promediocas = registrossem[contaminante].mean()
45 360     dataframesca[contaminante][i] = promediocas

```

Código 4.1: Procesamiento y agrupamiento de datos.

```

1 361 import pandas as pd
2 362 from epiweeks import Week, date
3 363 from sklearn import preprocessing
4 364 import seaborn as sns
5 365 import matplotlib.pyplot as plt
6 366 import string
7 367
8 368 diagnosticoaño = dataframea['DIAG_INI'].value_counts()
9 369 diagnosticoaño = diagnosticoaño.sort_values(ascending = False)
10 370 ciesaño = dataframea.groupby(['DIAG_INI', 'sem']).count()
11 371
12 372 s_scaler = preprocessing.StandardScaler()
13 373 ind = []
14 374 n = len(semanas.index)
15 375 for i in range(n):
16 376     ind.append(i+1)
17 377 letras = []
18 378 for letra in string.ascii_uppercase:
19 379     letras.append(str(letra))
20 380 # Se inicia un contador para controlar la cantidad de graficos a generar
21 381 cont = 0
22 382 maximo = 10
23 383 mindividuales = 7
24 384
25 385 # Proceso de generación de las figuras
26 386 print('\n' + año)
27 387 for name in diagnosticoaño.index:
28 388     if cont < maximo:
29 389         dataframegraficoacc = pd.DataFrame()
30 390         dataframegraficoacc[contaminante] = dataframesca[contaminante]
31 391         dataframegraficoacc = dataframegacc.reindex(ind)
32 392         if cont < mindividuales:
33 393             dataframegacc[name] = ciesaño['EGRESO'][name]
34 394             for i in range(n):
35 395                 dataframegacc[contaminante][i+1] = dataframesca[contaminante][i]
36 396             col_names = [contaminante, name]
37 397         else:
38 398             nameg = letras[cont]
39 399             ciesagrupadas = dataframea.loc[dataframea['DIAG_INI'].str.startswith(nameg)]
40 400             ciesagrupadas = ciesagrupadas['sem'].value_counts()
41 401             dataframegacc[nameg] = ciesagrupadas
42 402             for i in range(n):
43 403                 dataframegacc[contaminante][i+1] = dataframesca[contaminante][i]
44 404             col_names = [contaminante, nameg]
45 405         df_s = s_scaler.fit_transform(dataframegacc)
46 406         df_s = pd.DataFrame(df_s, columns=col_names)
47 407         fig, ax = plt.subplots(ncols=1, figsize=(20, 8))
48 408         print('\n' + col_names[0] + ' & ' + col_names[1])
49 409         ax.set_title('Contaminante ' + col_names[0] + ' & CIE ' + col_names[1])
50 410         ax.set_xlabel('Semana del año ' + año)
51 411         sns.kdeplot(data=df_s)
52 412         plt.savefig(contaminante + '/' + col_names[0] + '&' + col_names[1] + '_' + año + '.jpg',
413 format='jpg')
53 414         plt.show()
54 415         cont = cont+1

```

Código 4.2: Generación de series de tiempo.

```

1 416 def create_spiderwebs(datasets, labels, lenlines, title, titles, spoke_labels, colors, typeframe,
417     outputtype):
2 418
3 419     N1 = len(datasets)
4 420     N2 = len(labels)
5 421     N = int(N1/N2)
6 422     theta = radar_factory(N, frame=typeframe)
7 423     i=0
8 424     filenames = []
9 425     for titlespiderweb in titles:
10 426         fig, axs = plt.subplots(figsize=(8, 8), subplot_kw=dict(projection='radar'))
11 427         fig.subplots_adjust(wspace=0.5, hspace=0.20, top=0.85, bottom=0.05)
12 428         ax = axs
13 429         ax.set_title(titlespiderweb, weight='bold', size='medium', position=(0.5, 0.5),
430             horizontalalignment='center', verticalalignment='center', fontsize=16)
14 431         for y in range(N2):
15 432             dataspider = []
16 433             xx = y
17 434             for yy in range(N):
18 435                 currentdata = datasets[xx]
19 436                 number = currentdata[i]
20 437                 nmin = min(currentdata);
21 438                 nmax = max(currentdata);
22 439                 r = nmax - nmin
23 440                 x = (number-nmin)/r
24 441                 yyy = lenlines*x
25 442                 dataspider.append(yyy)
26 443                 xx = xx + N2
27 444                 ax.plot(theta, dataspider, color=colors[y])
28 445                 ax.fill(theta, dataspider, facecolor=colors[y], alpha=0.05)
29 446             ax.set_varlabels(spoke_labels)
30 447             ax = axs
31 448             legend = ax.legend(labels, loc=(0.9, .95), labelspacing=0.1, fontsize='small')
32 449             i=i+1
33 450             filename = 'spiderweb' + '_' + title + '_' + str(i) + '.jpg'
34 451             filenames.append(filename)
35 452             plt.savefig(filename, format='jpg')
36 453
37 454     # Generate a GIF
38 455     if(outputtype == 'gif'):
39 456         with imageio.get_writer(title + '.gif', mode='I', duration=1) as writer:
40 457             for filename in filenames:
41 458                 image = imageio.imread(filename)
42 459                 writer.append_data(image)
43 460     # Generate a .mp4 video
44 461     if(outputtype == 'video'):
45 462         img_array = []
46 463         for filename in filenames:
47 464             img = cv2.imread(filename)
48 465             height, width, layers = img.shape
49 466             size = (width,height)
50 467             img_array.append(img)
51 468         out = cv2.VideoWriter(title + '.mp4',cv2.VideoWriter_fourcc(*'MP4V'), 1, size)
52 469         for i in range(len(img_array)):
53 470             out.write(img_array[i])
54 471         out.release()

```

Código 4.3: Generación de graficos de radar.

```

1 472 # Gráfico
2 473 fig, ax = plt.subplots(figsize=(6, 3.84))
3 474 datos.plot(
4 475     x = col_names[0],
5 476     y = col_names[1],
6 477     c = 'firebrick',
7 478     kind = "scatter",
8 479     ax = ax
9 480 )
10 481 ax.set_title('Contaminante ' + col_names[0] + ' & CIE ' + col_names[1])
11 482 # Correlación lineal entre las dos variables
12 483 corr_test = pearsonr(x = datos[col_names[0]], y = datos[col_names[1]])
13 484 # División de los datos en train y test
14 485 X = datos[[col_names[0]]]
15 486 y = datos[col_names[1]]
16 487 X_train, X_test, y_train, y_test = train_test_split(
17 488     X.values.reshape(-1,1),
18 489     y.values.reshape(-1,1),
19 490     train_size = 0.8,
20 491     random_state = 1234,
21 492     shuffle = True
22 493 )
23 494 X_train = sm.add_constant(X_train, prepend=True)
24 495 modelo = sm.OLS(endog=y_train, exog=X_train)
25 496 modelo = modelo.fit()
26 497 if(modelo.pvalues[1]>0.05):
27 498     exclude_p.append(col_names[1])
28 499 # Intervalos de confianza para los coeficientes del modelo
29 500 modelo.conf_int(alpha=0.05)
30 501 # Predicciones con intervalo de confianza del 95%
31 502 pred = modelo.get_prediction(exog = X_train).summary_frame(alpha=0.05)
32 503 pred.head(4)
33 504 # Predicciones con intervalo de confianza del 95%
34 505 pred = modelo.get_prediction(exog = X_train).summary_frame(alpha=0.05)
35 506 pred['x'] = X_train[:, 1]
36 507 pred['y'] = y_train
37 508 pred = pred.sort_values('x')
38 509 # Gráfico del modelo
39 510 fig, ax = plt.subplots(figsize=(6, 3.84))
40 511 ax.set_title('Contaminante ' + col_names[0] + ' & CIE ' + col_names[1])
41 512 ax.set_xlabel(col_names[0])
42 513 ax.set_ylabel(col_names[1])
43 514 ax.scatter(pred['x'], pred['y'], marker='o', color = "gray")
44 515 ax.plot(pred['x'], pred["mean"], linestyle='-', label="OLS")
45 516 ax.plot(pred['x'], pred["mean_ci_lower"], linestyle='--', color='red')
46 517 ax.plot(pred['x'], pred["mean_ci_upper"], linestyle='--', color='red')
47 518 ax.fill_between(pred['x'], pred["mean_ci_lower"], pred["mean_ci_upper"], 0.1)
48 519 ax.legend()
49 520 # Error de test del modelo
50 521 X_test = sm.add_constant(X_test, prepend=True)
51 522 pred = modelo.predict(exog = X_test)
52 523 rmse = mean_squared_error(
53 524     y_true = y_test,
54 525     y_pred = pred,
55 526     squared = False
56 527 )

```

Código 4.4: Generación de los modelos de regresión lineal.

```

1 528 corr_matrix = datarlm.corr(method='pearson')
2 529 corr_mat = corr_matrix.stack().reset_index()
3 530 corr_mat.columns = ['variable_1', 'variable_2', 'r']
4 531 corr_mat = corr_mat.loc[corr_mat['variable_1'] != corr_mat['variable_2'], :]
5 532 corr_mat['abs_r'] = np.abs(corr_mat['r'])
6 533 corr_mat = corr_mat.sort_values('abs_r', ascending=False)
7 534 tidy_corr_matrix = (corr_matrix).head(10)
8 535 # Heatmap matriz de correlaciones
9 536 fig, ax = plt.subplots(nrows=1, ncols=1, figsize=(4, 4))
10 537 sns.heatmap(
11 538     corr_matrix,
12 539     annot      = True,
13 540     cbar       = False,
14 541     annot_kws  = {"size": 8},
15 542     vmin       = -1,
16 543     vmax       = 1,
17 544     center     = 0,
18 545     cmap       = sns.diverging_palette(20, 220, n=200),
19 546     square     = True,
20 547     ax         = ax)
21 548 ax.set_xticklabels(
22 549     ax.get_xticklabels(),
23 550     rotation = 45,
24 551     horizontalalignment = 'right',)
25 552 ax.tick_params(labelsize = 10)
26 553 # División de los datos en train y test
27 554 X = datarlm[spoke_labels]
28 555 y = datarlm[contaminante]
29 556 X_train, X_test, y_train, y_test = train_test_split(
30 557     X,
31 558     y.values.reshape(-1,1),
32 559     train_size    = 0.8,
33 560     random_state  = 1234,
34 561     shuffle       = True)
35 562 # Creación del modelo utilizando matrices como en scikitlearn
36 563 X_train = sm.add_constant(X_train, prepend=True)
37 564 modelo = sm.OLS(endog=y_train, exog=X_train,)
38 565 modelo = modelo.fit()
39 566 sml = modelo.summary().as_latex()
40 567 namefile = 'modelos_latex/' + 'regresion_lineal_multiple_' + contaminante + '_' + año + '.tex'
41 568 f = open(namefile, 'w')
42 569 with open(namefile, 'w') as f:
43 570     f.write(sml)
44 571 # Diagnóstico errores (residuos) de las predicciones de entrenamiento
45 572 y_train = y_train.flatten()
46 573 predicciones_train = modelo.predict(exog = X_train)
47 574 residuos_train = predicciones_train - y_train
48 575 # Predicciones con intervalo de confianza
49 576 predicciones = modelo.get_prediction(exog = X_train).summary_frame(alpha=0.05)
50 577 predicciones.head(4)
51 578 # Error de test del modelo
52 579 X_test = sm.add_constant(X_test, prepend=True)
53 580 rmse = mean_squared_error(
54 581     y_true = y_test,
55 582     y_pred = modelo.predict(exog = X_test),
56 583     squared = False)

```

Código 4.5: Generación de los modelos de regresión lineal múltiple.

584

CAPÍTULO 5

585

EXPERIMENTOS

586

En el presente capítulo se presenta el diseño de los experimentos realizados así como los resultados obtenidos de ellos.

588

En esta sección se tratan los resultados obtenidos partiendo de desarrollar algunos experimentos que permiten determinar si la solución propuesta cumple con el objetivo planteado.

591

5.1 DISEÑO EXPERIMENTAL

592

En la presente sección se discute el diseño de experimentos, es decir, qué valores constantes fueron utilizados para su realización y por qué se usan dichos valores.

594

5.1.1 DATOS DE ENTRADA

595

Los datos de egresos hospitalarios de los años 2017 y 2018 provienen de la base de datos de la Secretaría de Salud del Gobierno de México [9]. También se tienen registros de los niveles de PM10 y PM2.5 presentes en el área metropolitana de Monterrey, dichos registros son obtenidos por las estaciones de monitoreo pertenecientes

599 al SIMA [20] mostradas en la figura 1.1 en la página 2. Los documentos con los datos
600 son proporcionados por Benavides [3].

601 **Selección de datos.** Los conjuntos de datos por año de egresos hospitalarios con-
602 tienen información de todos los estados de México, por lo cual se hace una
603 limpieza de datos para solo obtener los registros de Nuevo León ya que de
604 dicha entidad es de la cual se tienen los datos de contaminación.

605 **Datos de ingresos hospitalarios.** Se agrupan en semanas epidemiológicas para
606 una mejor manipulación de ellos. Por CIE se obtiene el número de egresos en
607 cada semana del año ordenadas por el número de egresos en el año partiendo
608 de la CIE con la mayor cantidad.

609 **Datos de los contaminantes.** Se agrupan en semanas epidemiológicas para una
610 mejor manipulación de ellos. Se obtiene el promedio del nivel del contaminante
611 por cada semana del año.

612 5.1.2 VISUALIZACIÓN DE DATOS

613 Con los datos ya seleccionados y agrupados se procede a generar las visualiza-
614 ciones de los datos. Las visualizaciones de los datos se hacen con series de tiempo y
615 gráficos de radar.

616 **Series de tiempo.** Se procede a generar series de tiempo de cada año por conta-
617 minante y CIE. Las variables ajustables son:

- 618 ■ El nombre del contaminante.
- 619 ■ El año del que se quieren obtener las series de tiempo.
- 620 ■ Número de series de tiempo a generar por contaminante. Se parte de la
621 CIE con mayor número de egresos.

622 **Gráficos de radar.** Los datos ya seleccionados y agrupados se normalizan teniendo
623 como valor mínimo cero y como valor máximo un número entre uno y cuatro.
624 Posteriormente se generan gráficos de radar de cada año por semana, en las
625 que se muestra el nivel contaminante y la variación de las CIE. Las variables
626 ajustables son:

- 627 ■ Cantidad de CIE a agregar en el gráfico. Se parte de la CIE con mayor
628 número de egresos.
- 629 ■ Valor máximo que se utiliza para representar la longitud de los ejes en el
630 gráfico.
- 631 ■ Nombre de la figura.
- 632 ■ Nombre de cada eje en el gráfico.
- 633 ■ Los colores de cada eje en el gráfico.
- 634 ■ Si el gráfico es generado de forma circular o en forma de polígono.

635 5.1.3 GENERACIÓN DE MODELOS

636 Se procede a generar los modelos. En cada modelo se tienen métricas para
637 evaluar su eficacia y valores que pueden ser ajustados en función de encontrar la
638 combinación que proporcione mejores resultados.

639 Se tienen algunas variables que pueden ser modificadas para la generación de
640 los modelos de regresión lineal.

- 641 ■ Cantidad de CIE a agregar en el modelo. Se parte de la CIE con mayor número
642 de egresos.
- 643 ■ Porcentaje de datos utilizados para el entrenamiento del modelo.
- 644 ■ Nivel de significancia.

645 También se tienen variables que indican información sobre la eficacia del mo-
646 delo.

647 ■ Valor p .

648 ■ R^2 (R cuadrado).

649 ■ Raíz de error cuadrático medio (RMSE).

5.2 RESULTADOS

Establecidas las especificaciones de los experimentos que se realizan, se reportan los resultados obtenidos. Los experimentos se elaboran por contaminante, desglosando los resultados por año. En la carpeta <https://github.com/selenebpradop/relaciones-contaminantes-salud/tree/main/figuras/> se encuentran animaciones en video de los gráficos de radar generados por contaminante y año y todas las imágenes de las series de tiempo obtenidas.

5.2.1 EXPERIMENTO A: DATOS DE NIVELES DE PM10

Se estudian los niveles del contaminante PM10 de los años 2017 y 2018.

5.2.1.1 AÑO 2017

En la figura 5.1 se muestra una de las series de tiempo generadas para la CIE con mayor número de egresos registrados en el conjunto de datos del año. Además, en el cuadro 5.1 se presentan los resultados obtenidos de los modelos de regresión lineal y la eficacia obtenida de dichos modelos. El cuadro 5.2 muestra los resultados del modelo de regresión lineal múltiple. En la figura 5.2 se muestran las correlaciones obtenidas en un gráfico de radar.

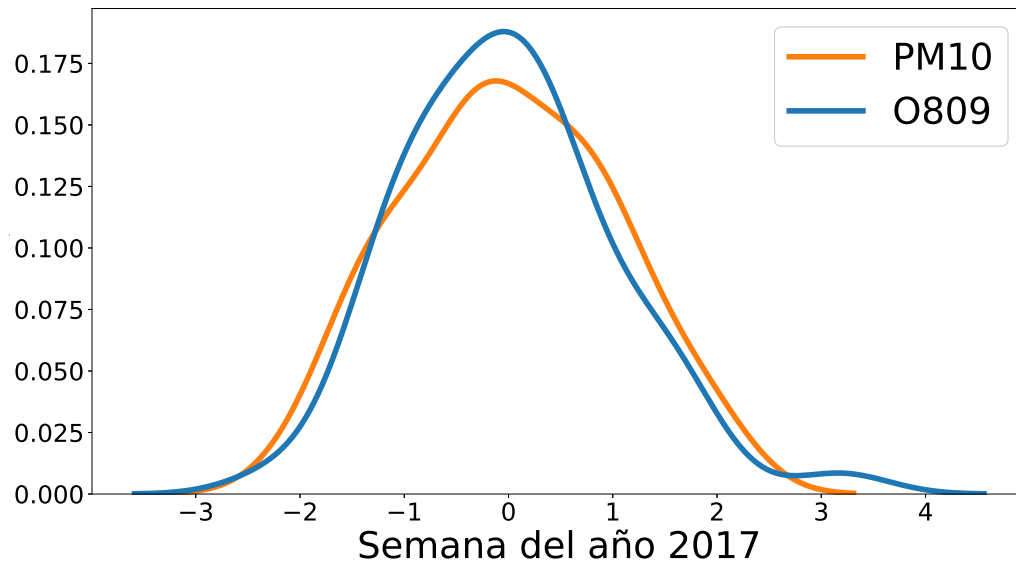


Figura 5.1: Evolución de los niveles de PM10 y el número de egresos diagnosticados con la CIE O809 en el 2017.

Cuadro 5.1: Resultados obtenidos PM10 2017

CIE	ρ	R^2	Valor p	ϵ
O809	-0.275	0.061	0.121	0.239
O829	0.100	0.091	0.055	0.287
O759	-0.085	0.116	0.029	0.294
O069	-0.247	0.070	0.094	0.222
K802	0.044	0.005	0.658	0.282

ρ = Coeficiente de correlación de Pearson

ϵ = RMSE para medir el error

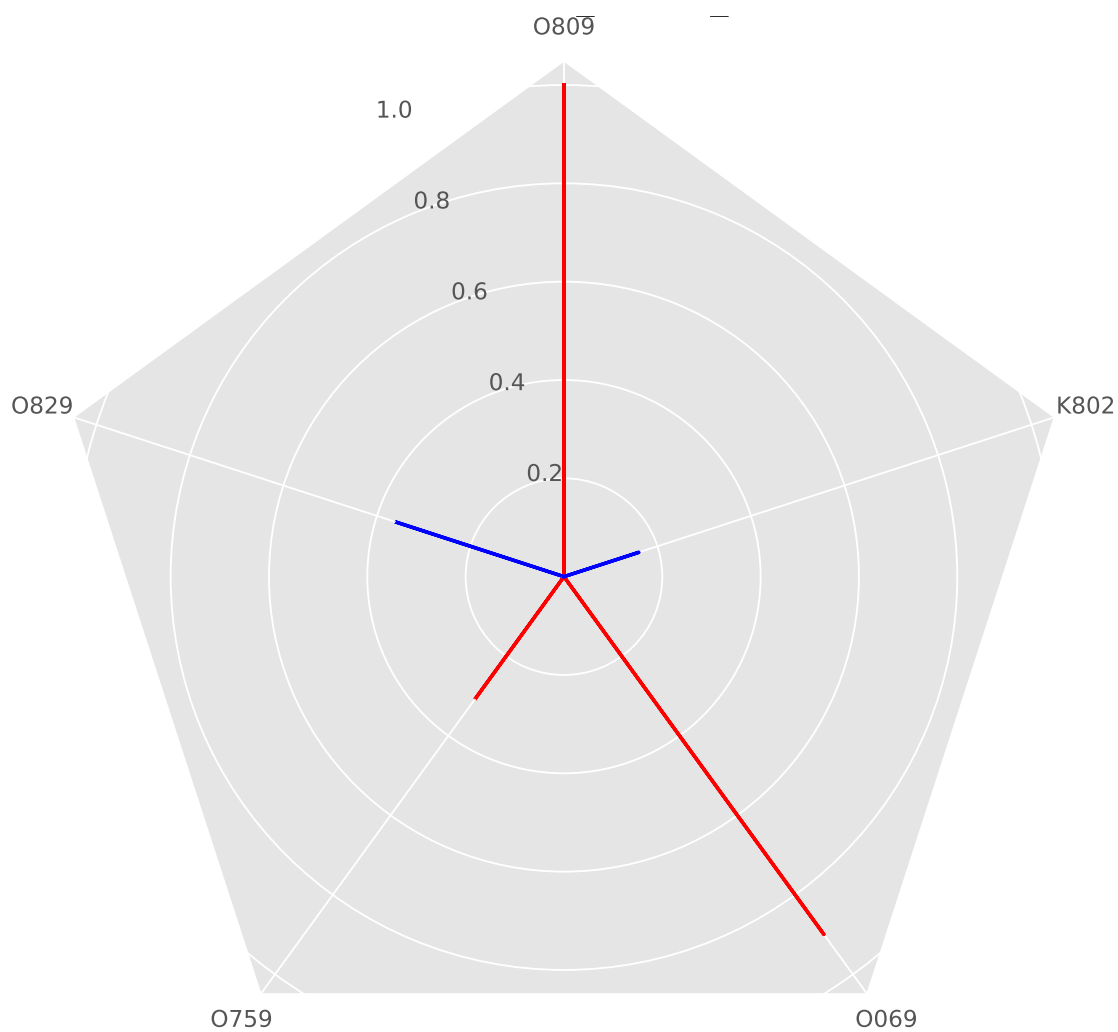


Figura 5.2: Correlaciones entre los niveles de PM10 y CIE en el 2017 donde el azul indica una correlación positiva y el rojo una correlación negativa.

Cuadro 5.2: Resultados regresión lineal múltiple PM10 2017

Variable Dep.:	y	R²:	0.313			
Modelo:	OLS	Método:	Mínimos cuadrados			
Error:	0.226					
	coef	std err	t	P> t	[0.025	0.975]
const	0.5859	0.146	4.002	0.000	0.289	0.883
O809	-0.3019	0.122	-2.472	0.018	-0.550	-0.054
O829	0.2685	0.123	2.185	0.036	0.019	0.518
O759	-0.2229	0.182	-1.222	0.230	-0.593	0.147
O069	-0.1441	0.120	-1.200	0.238	-0.388	0.100
K802	-0.0456	0.133	-0.344	0.733	-0.315	0.224

5.2.1.2 Año 2018

En la figura 5.3 se muestra una de las series de tiempo generadas para la CIE con mayor número de egresos registrados en el conjunto de datos del año. Además, en el cuadro 5.3 se presentan los resultados obtenidos de los modelos de regresión lineal y la eficacia obtenida de dichos modelos. El cuadro 5.4 muestra los resultados del modelo de regresión lineal múltiple. En la figura 5.4 se muestran las correlaciones obtenidas en un gráfico de radar.

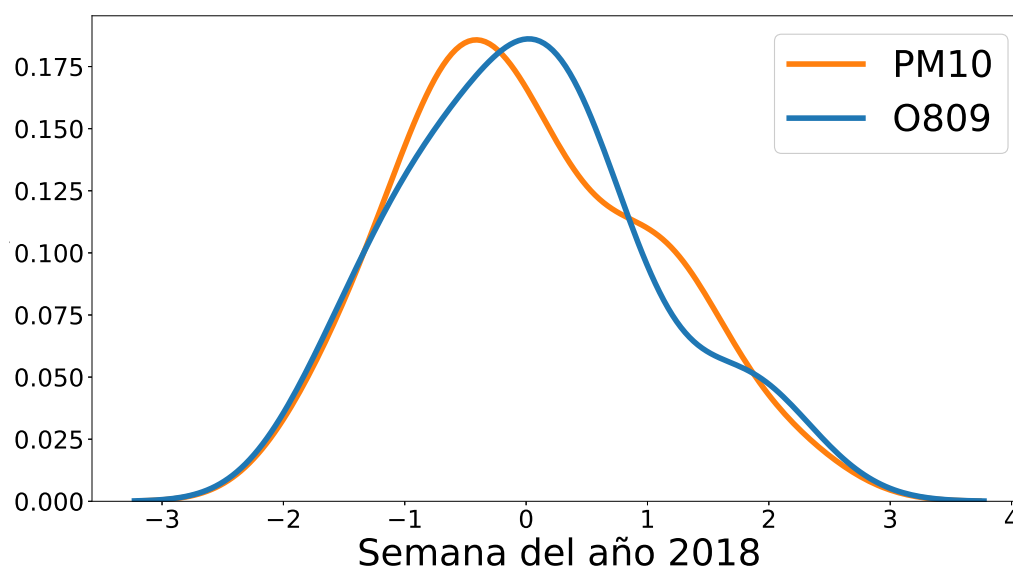


Figura 5.3: Evolución de los niveles de PM10 y el número de egresos diagnosticados con la CIE O809 en el 2018.

Cuadro 5.3: Resultados obtenidos PM10 2018

CIE	ρ	R^2	Valor p	ϵ
O809	-0.271	0.015	0.440	0.275
O829	0.282	0.069	0.098	0.249
K802	0.277	0.084	0.066	0.152
O342	-0.401	0.100	0.044	0.248
N40X	-0.009	0.000	0.964	0.243

ρ = Coeficiente de correlación de Pearson

ϵ = RMSE para medir el error

Cuadro 5.4: Resultados regresión lineal múltiple PM10 2018

Variable Dep.:	y	R²:	0.182			
Modelo:	OLS	Método:	Mínimos cuadrados			
Error:	0.159					
	coef	std err	t	P> t	[0.025	0.975]
const	0.3885	0.210	1.849	0.073	-0.038	0.815
O809	-0.0114	0.188	-0.061	0.952	-0.392	0.370
O829	0.0854	0.214	0.400	0.692	-0.349	0.520
K802	0.3088	0.167	1.852	0.072	-0.030	0.647
O342	-0.1708	0.208	-0.820	0.418	-0.593	0.252
N40X	-0.1239	0.167	-0.740	0.464	-0.464	0.216

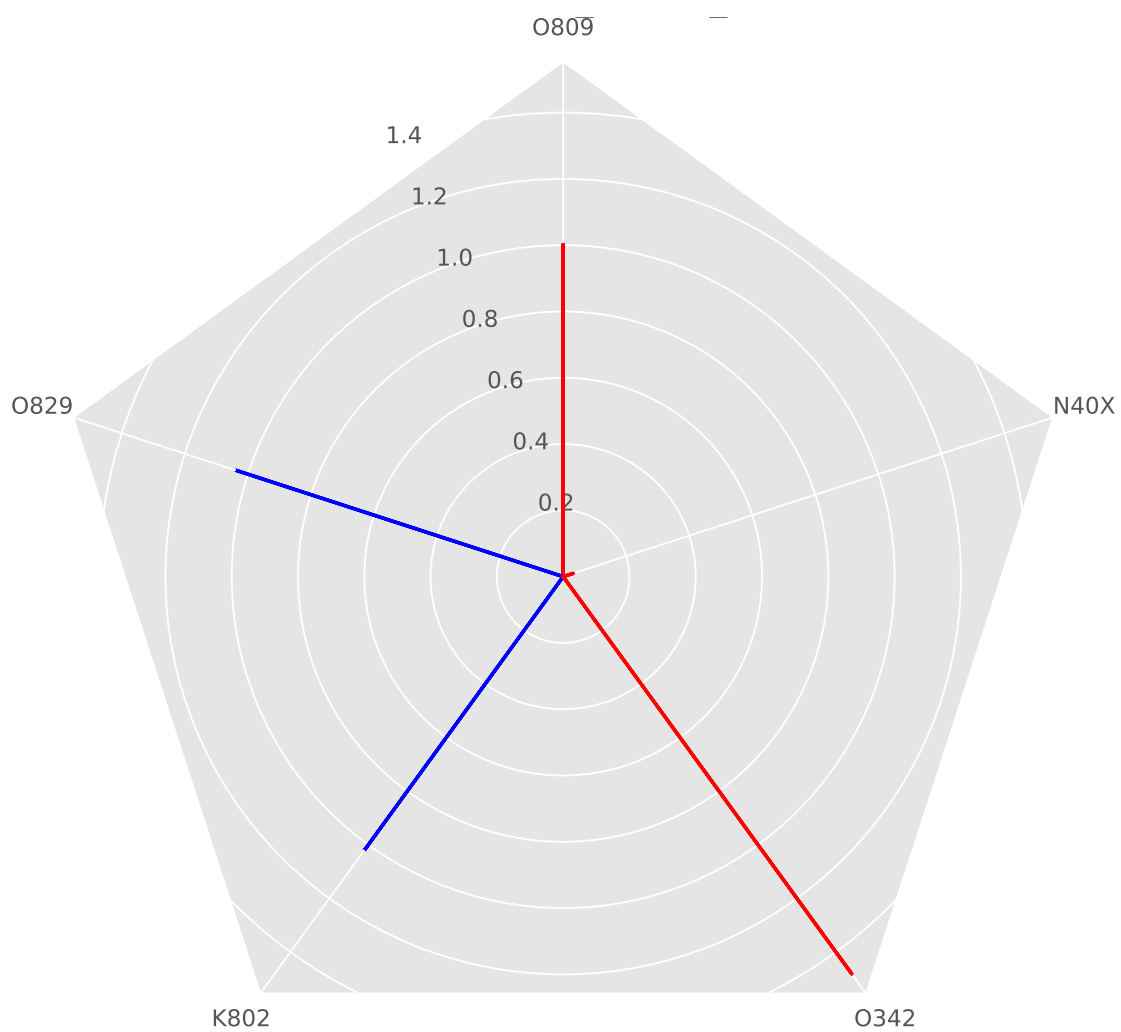


Figura 5.4: Correlaciones entre los niveles de PM10 y CIE en el 2018 donde el azul indica una correlación positiva y el rojo una correlación negativa.

5.2.2 EXPERIMENTO B: NIVELES DE PM2.5

Se estudian los niveles del contaminante PM2.5 de los años 2017 y 2018.

5.2.2.1 AÑO 2017

En la figura 5.5 se muestra una de las series de tiempo generadas para la CIE con mayor número de egresos registrados en el conjunto de datos del año. Además, en el cuadro 5.5 se presentan los resultados obtenidos de los modelos de regresión lineal y la eficacia obtenida de dichos modelos. El cuadro 5.6 muestra los resultados del modelo de regresión lineal múltiple. En la figura 5.6 se muestran las correlaciones obtenidas en un gráfico de radar.

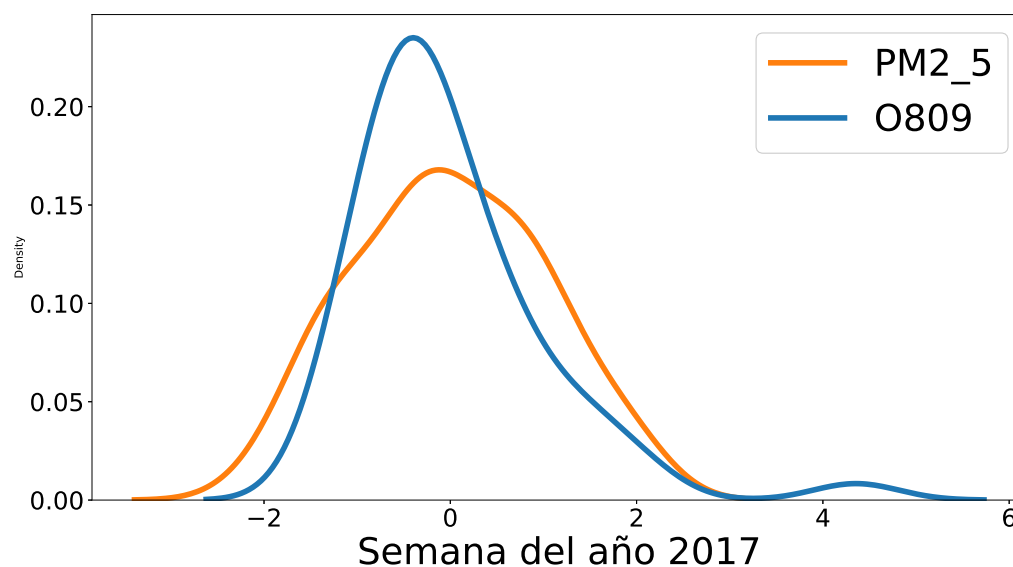


Figura 5.5: Evolución de los niveles de PM2.5 y el número de egresos diagnosticados con la CIE O809 en el 2017.

Cuadro 5.5: Resultados obtenidos PM2.5 2017

CIE	ρ	R^2	Valor p	ϵ
O809	-0.093	0.011	0.511	0.256
O829	0.014	0.021	0.371	0.253
O759	-0.172	0.113	0.032	0.278
O069	-0.350	0.112	0.033	0.208
K802	0.006	0.005	0.667	0.285

ρ = Coeficiente de correlación de Pearson

ϵ = RMSE para medir el error

Cuadro 5.6: Resultados regresión lineal múltiple PM2.5 2017

Variable Dep.:	y	R²:	0.210			
Modelo:	OLS	Método:	Mínimos cuadrados			
Error:	0.175					
	coef	std err	t	P> t	[0.025	0.975]
const	0.5345	0.158	3.373	0.002	0.213	0.856
O809	-0.1754	0.132	-1.327	0.193	-0.444	0.093
O829	0.0701	0.133	0.527	0.602	-0.200	0.340
O759	-0.2585	0.197	-1.309	0.199	-0.659	0.142
O069	-0.2148	0.130	-1.652	0.108	-0.479	0.049
K802	-0.1164	0.144	-0.811	0.423	-0.408	0.175

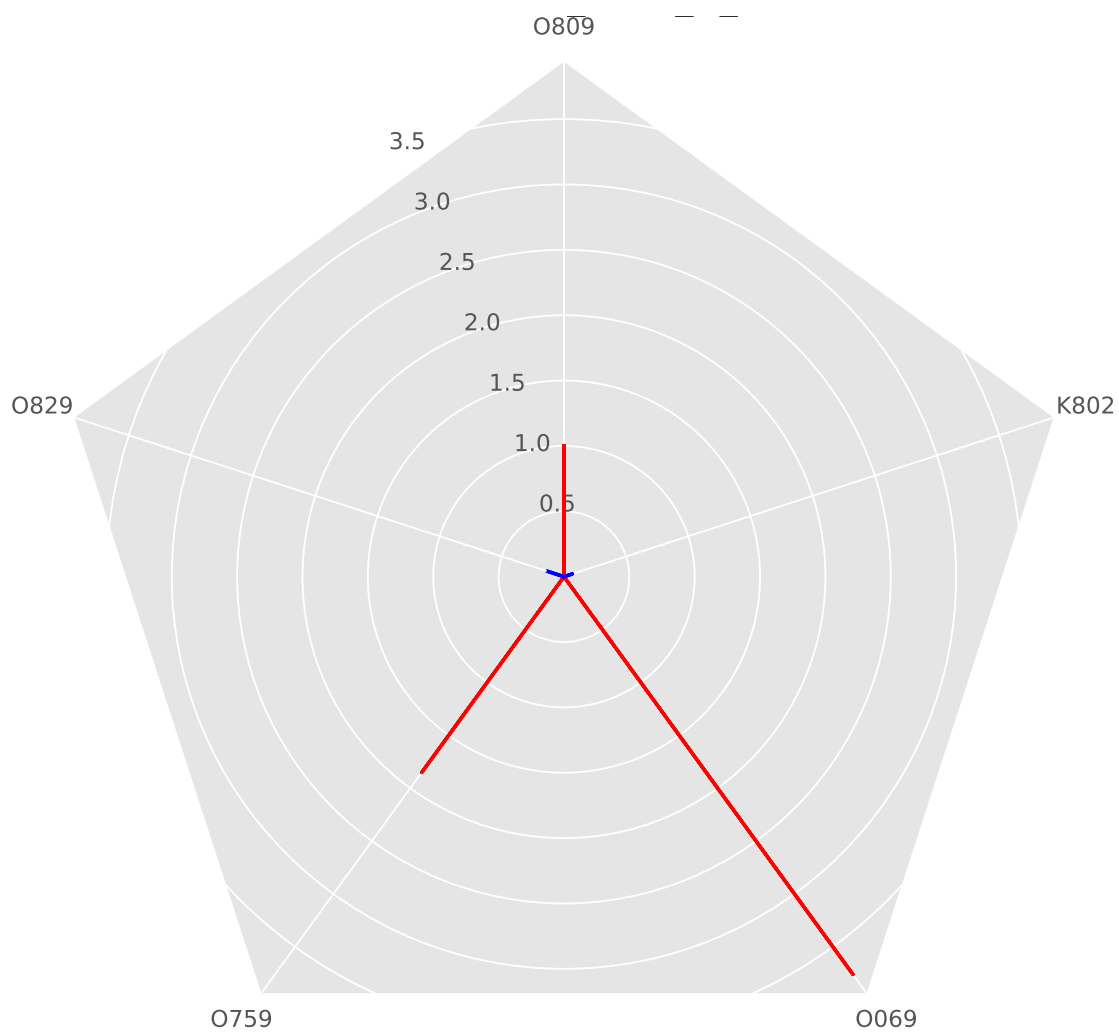


Figura 5.6: Correlaciones entre los niveles de PM2.5 y CIE en el 2017 donde el azul indica una correlación positiva y el rojo una correlación negativa.

5.2.2.2 Año 2018

En la figura 5.7 se muestra una de las series de tiempo generadas para la CIE con mayor número de egresos registrados en el conjunto de datos del año. Además, en el cuadro 5.7 se presentan los resultados obtenidos de los modelos de regresión lineal y la eficacia obtenida de dichos modelos. El cuadro 5.8 muestra los resultados del modelo de regresión lineal múltiple. En la figura 5.8 se muestran las correlaciones obtenidas en un gráfico de radar.

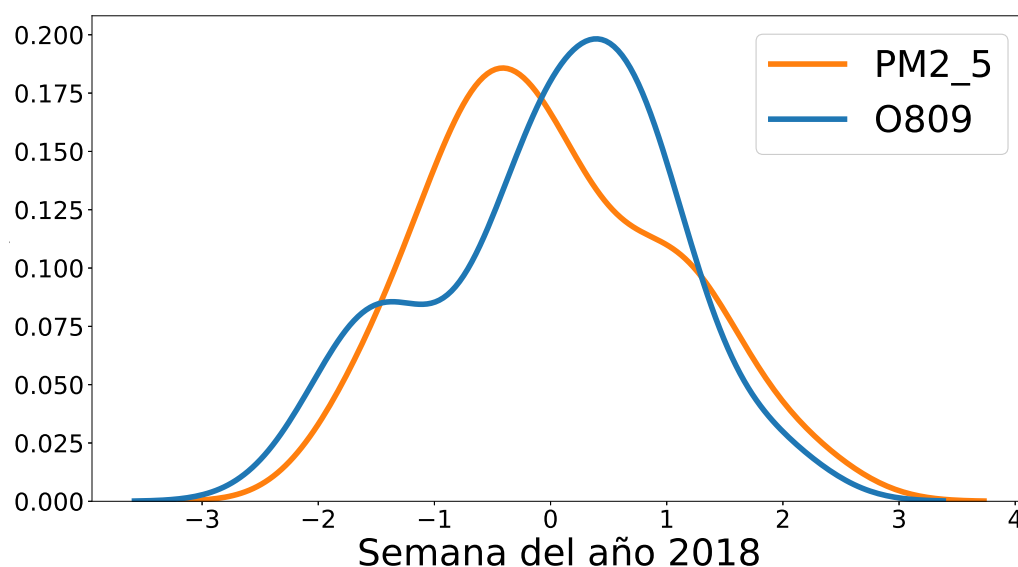


Figura 5.7: Evolución de los niveles de PM2.5 y el número de egresos diagnosticados con la CIE O809 en el 2018.

Cuadro 5.7: Resultados obtenidos PM2.5 2018

CIE	ρ	R^2	Valor p	ϵ
O809	-0.354	0.057	0.133	0.259
O829	0.225	0.046	0.177	0.256
K802	0.273	0.089	0.058	0.159
O342	-0.443	0.149	0.013	0.245
N40X	-0.074	0.001	0.842	0.241

ρ = Coeficiente de correlación de Pearson

ϵ = RMSE para medir el error

Cuadro 5.8: Resultados regresión lineal múltiple PM2.5 2018

Variable Dep.:	y	R²:	0.259			
Modelo:	OLS	Método:	Mínimos cuadrados			
Error:	0.174					
	coef	std err	t	P> t	[0.025	0.975]
const	0.7042	0.193	3.652	0.001	0.313	1.096
O809	-0.0863	0.172	-0.501	0.620	-0.436	0.264
O829	-0.1084	0.196	-0.553	0.584	-0.507	0.290
K802	0.3006	0.153	1.964	0.058	-0.010	0.611
O342	-0.3311	0.191	-1.733	0.092	-0.719	0.057
N40X	-0.2053	0.154	-1.336	0.190	-0.517	0.107

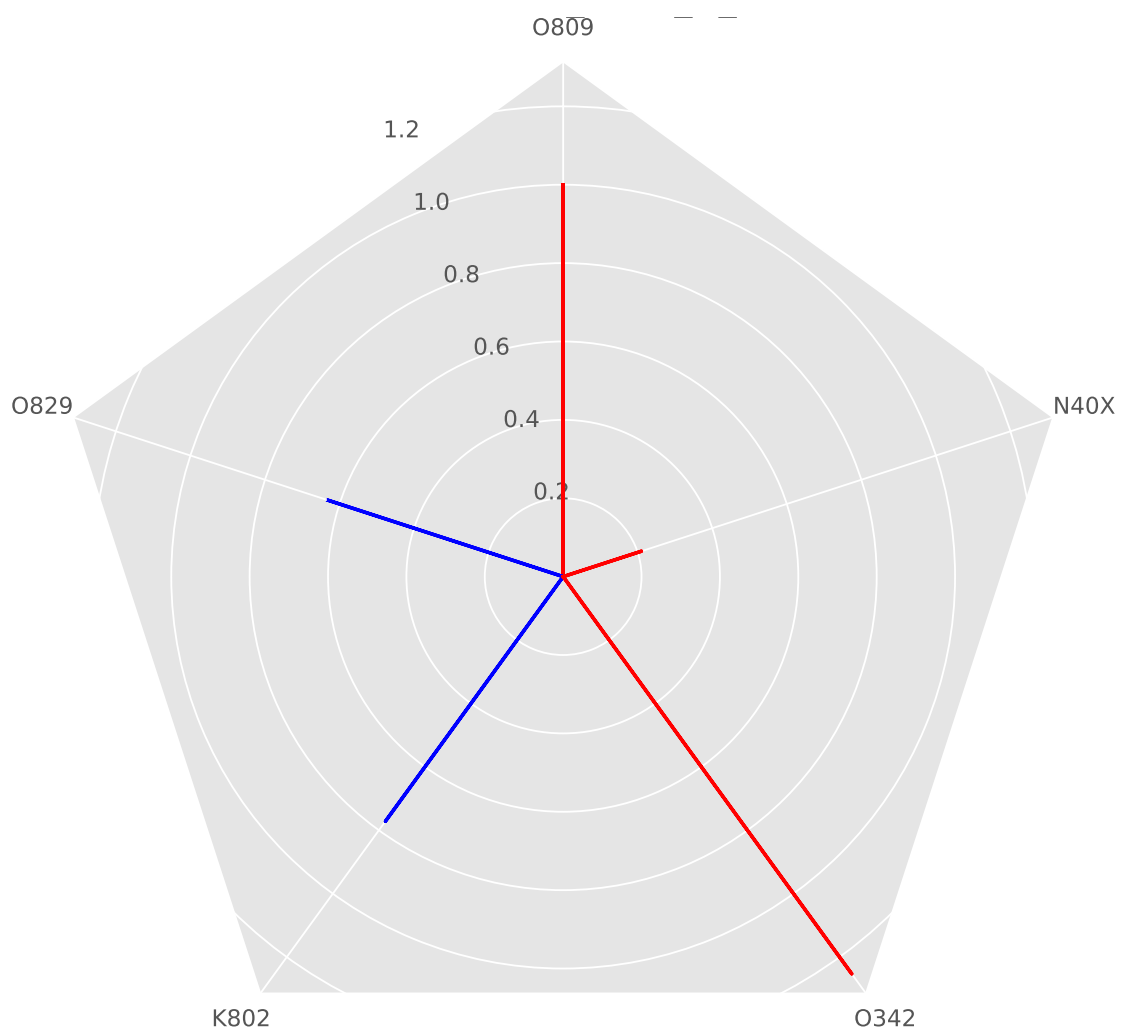


Figura 5.8: Correlaciones entre los niveles de PM2.5 y CIE en el 2018 donde el azul indica una correlación positiva y el rojo una correlación negativa.

5.3 DISCUSIÓN

689 Como se puede observar, en todos los experimentos se obtiene una correlación
690 entre -1 y 1 y diferente a 0, por lo tanto se pueden generar los modelos de regresión
691 lineal.
692

693 En el Experimento A el error RMSE en los modelos de regresión lineal varía
694 entre 0.152 y 0.294, lo cual indica que no todos los resultados alcanzan una fiabilidad
695 mayor al 80 %, excepto en la CIE K802 en el año 2018 en la cual se encuentra el
696 porcentaje de error más bajo. El valor de R^2 más alto en el año 2017 se encuentra en
697 la CIE O759 con un valor p de 0.029, sin embargo, en el modelo de regresión lineal
698 múltiple se encuentra el valor de R^2 más alto para el contaminante PM10 en el año
699 2017. En el año 2018 el valor de R^2 más alto se encuentra en la CIE O342 con un
700 valor p de 0.044, sin embargo, en el modelo de regresión lineal múltiple se encuentra
701 el valor de R^2 más alto para el contaminante PM10 en el año 2018.

702 En el Experimento B el error RMSE en los modelos de regresión lineal varía
703 entre 0.159 y 0.278, lo cual indica que no todos los resultados alcanzan una fiabilidad
704 mayor al 80 %, excepto en la CIE K802 en el año 2018 en la cual se encuentra el
705 porcentaje de error más bajo. El valor de R^2 más alto en el año 2017 se encuentra en
706 la CIE O759 con un valor p de 0.032, sin embargo, en el modelo de regresión lineal
707 múltiple se encuentra el valor de R^2 más alto para el contaminante PM2.5 en el año
708 2017. En el año 2018 el valor de R^2 más alto se encuentra en la CIE O342 con un
709 valor p de 0.013, sin embargo, en el modelo de regresión lineal múltiple se encuentra
710 el valor de R^2 más alto para el contaminante PM2.5 en el año 2018.

711 Todos los experimentos son ejecutados en `Jupyter Notebook` en una laptop
712 con las especificaciones del cuadro 5.9.

Cuadro 5.9: Especificaciones técnicas del equipo de cómputo

Sistema Operativo	macOS Big Sur
Procesador	Apple M1
RAM	8 GB RAM

CONCLUSIONES

715 El presente capítulo describe la tesis a partir de la manera que cumple los
716 objetivos generales y específicos para determinar si la hipótesis se comprueba, trata
717 también del porqué se realizó la tesis.

718 En el presente proyecto se generaron visualizaciones para el estudio de las rela-
719 ciones entre determinados contaminantes y determinadas CIE. Además, se generaron
720 modelos de regresión lineal y modelos de regresión lineal múltiple con diferentes con-
721 taminantes y diferentes CIE para poder realizar una comparación entre los modelos
722 generados.

723 En los experimentos se encontró un coeficiente de correlación de Pearson di-
724 ferente a 0 entre -1 y 1, por lo tanto se pudieron generar los modelos de regresión
725 lineal. Se encontró que al estudiar las CIE de manera agrupada en un modelo de
726 regresión múltiple se obtiene una mejor explicación de la varianza de los niveles de
727 los contaminantes PM10 y PM2.5 frente al modelo de regresión lineal simple. Cinco
728 de los veinticuatro valores de error RMSE reportados son menores a 0.20, lo cual
729 indica que la mayoría de los valores predichos por los modelos se alejaron más de
730 0.20 unidades de los valores reales, por lo tanto dichos valores de error son mayores
731 a los deseados.

6.1 CONTRIBUCIONES

732

733 Primeramente se encontró que la CIE que reporta mayor número de egresos en
734 Nuevo León, México en los años 2017 y 2018 es la CIE O809. En las series de tiempo
735 se observa que la cantidad de egresos de la mayoría de las CIE estudiadas presentan
736 una línea de evolución similar al contaminante PM10.

737 Además, se encontró una correlación lineal entre los contaminantes estudiados
738 y las CIE estudiadas, por lo cual el presente proyecto motiva a seguir realizando
739 labores en torno a la investigación de la dependencia lineal que se presenta entre los
740 contaminantes y las CIE.

741 Finalmente, se observó que para el estudio de las relaciones entre los contami-
742 nantes y las CIE se obtuvieron mejores resultados con los modelos de regresión lineal
743 múltiple frente a los modelos de regresión lineal simple, lo cual indica que se puede
744 obtener información relevante si se emplean modelos de regresión lineal múltiple para
745 realizar investigaciones que estudien las relaciones entre los niveles de determinados
746 contaminantes y el número de egresos por determinadas CIE.

6.2 TRABAJO A FUTURO

747

748 El presente trabajo brinda algunos aspectos a considerar para realizar un traba-
749 jo a futuro, los cuales son: la recolección de más datos de los niveles de contaminantes
750 y de egresos para el estudio de años más recientes, realizar un estudio de las relacio-
751 nes de los niveles de determinados contaminantes y la cantidad de egresos por CIE
752 empleando modelos de regresión lineal múltiple, la creación de un mapa interactivo
753 donde se pueda observar por región los niveles de los contaminantes y la cantidad de
754 egresos, y la generación de una página web que funcione en conjunto con el desarrollo
755 elaborado para que al ingresar el archivo con los datos se generen las visualizaciones
756 y modelos.

757 Las oportunidades de mejora detectada para el presente proyecto son: comparar
758 los resultados obtenidos de los modelos regresión lineal múltiple con el mismo modelo
759 pero incrementando el número de variables independientes (CIE), elaborar modelos
760 de regresión no lineal para comparar sus resultados con los resultados de los modelos
761 generados, y utilizar un conjunto de datos más consistente ya que eso favorece a que
762 se tengan resultados más fiables y precisos.

APÉNDICE A

CIE Y SUS NOMBRES DE ENFERMEDADES

Cuadro A.1: CIE mencionadas en los Experimentos y el nombre de la enfermedad.

CIE	Enfermedad
N40-N53	Enfermedades de los órganos genitales masculino
K00-K93	Enfermedades del aparato digestivo
O00-O99	Embarazo, parto y puerperio

BIBLIOGRAFÍA

- [1] ARIAS, J. R. (2006), «What is an epidemiological week and why do we use them», *Skeeter*, **66**(1), pág. 7.
- [2] BALLESTER DÍEZ, F., J. M. TENÍAS y S. PÉREZ-HOYOS (1999), «Efectos de la contaminación atmosférica sobre la salud: una introducción», *Revista Española de Salud Pública*, **73**(2), págs. 109–121.
- [3] BENAVIDES, A. (2022), «Perfil de Github», <https://github.com/jbenavidesv87>.
- [4] BRETON, R. M. C., J. C. BRETON, M. DE LA LUZ ESPINOSA FUENTES, J. KAHL, A. A. E. GUZMAN, R. G. MARTÍNEZ, C. GUARNACCIA, R. DEL CARMEN LARA SEVERINO, E. R. LARA y A. B. FRANCAVILLA (2021), «Short-Term Associations between Morbidity and Air Pollution in Metropolitan Area of Monterrey, Mexico», *Atmosphere*, **12**(10), pág. 1352. doi: 10.3390/atmos12101352.
- [5] BROCKWELL, P. J., P. J. BROCKWELL, R. A. DAVIS y R. A. DAVIS (2002), *Introduction to time series and forecasting*, Springer. ISBN 9780387216577.
- [6] CATALÁ, F. y E. DE MANUEL (1998), «Informe SESPAS 1998: La salud pública y el futuro del estado del bienestar», *Granada: EASP*.
- [7] CORBITT, R. y R. CORBITT (1990), *Standard Handbook of Environmental Engineering*, McGraw-Hill. ISBN 9780070131583.

-
- [8] DARLINGTON, R. B. y A. F. HAYES (2016), *Regression analysis and linear models: Concepts, applications, and implementation*, Guilford Publications. ISBN 9781462521135.
- [9] DIRECCIÓN GENERAL DE INFORMACIÓN EN SALUD (2021), «Egresos Hospitalarios», http://www.dgis.salud.gob.mx/contenidos/basesdedatos/da_egresoshosp_gobmx.html.
- [10] GUARNACCIA, C., J. G. C. BRETON, R. M. C. BRETON, C. TEPEDINO, J. QUARTIERI y N. E. MASTORAKIS (2018), «ARIMA models application to air pollution data in Monterrey, Mexico», en *AIP Conference Proceedings*, tomo 1, AIP Publishing LLC, págs. 1–5.
- [11] GUPTA, A., H. BHERWANI, S. GAUTAM, S. ANJUM, K. MUSUGU, N. KUMAR, A. ANSHUL y R. KUMAR (2021), «Air pollution aggravating COVID-19 lethality? Exploration in Asian cities using statistical models», *Environment, Development and Sustainability*, **23**(4), págs. 6408–6417. doi: 10.1007/s10668-020-00878-9.
- [12] HADDAD, R. (1974), «Contaminación del aire. Situación actual en la América Latina y el Caribe», *Informe técnico*.
- [13] JULIA, A., D. A. FC LICHTENFELS, K. VAN DER PLAAT, C. C. DE JONG, D. S. VAN DIEMEN, I. POSTMA, C. M. NEDELJKOVIC, N. VAN DUIJN, S. AMIN, M. LA BASTIDE-VAN GEMERT, D. VRIES *et al.* (2018), «Long-term air pollution exposure, genome-wide DNA methylation and lung function in the LifeLines cohort study», *Environmental health perspectives*, **126**(2), pág. 027004. doi: 10.1289/EHP2045.
- [14] KIM, J. S., Z. CHEN, T. L. ALDERETE, C. TOLEDO-CORRAL, F. LURMANN, K. BERHANE y F. D. GILLILAND (2019), «Associations of air pollution, obesity and cardiometabolic health in young adults: The Meta-AIR study», *Environment international*, **133**, págs. 105–180. doi: 10.1016/j.envint.2019.105180.

- [15] KORC, M. y R. SÁENZ (1999), «Monitoreo de la calidad del aire en América Latina», *Korc Marcelo E*, págs. 1–22.
- [16] LIU, Y., J. PAN, H. ZHANG, C. SHI, G. LI, Z. PENG, J. MA, Y. ZHOU y L. ZHANG (2019), «Short-term exposure to ambient air pollution and asthma mortality», *American journal of respiratory and critical care medicine*, **200**(1), págs. 24–32. doi: 10.1164/rccm.201810-1823OC.
- [17] MARTÍN, R. M. y M. S. BAYLE (2018), «Impacto de la contaminación ambiental en las consultas pediátricas de Atención Primaria: estudio ecológico», en *Anales de Pediatría*, tomo 2, Elsevier, págs. 80–85.
- [18] ORGANIZATION, W. H. *et al.* (2016), «International Classification of Diseases. 2016», *World Health Organization*.
- [19] PENICHE-CAMPS, S. y M. CORTEZ-HUERTA (2020), «La costumbre al envenenamiento: El caso de los contaminantes atmosféricos de la ciudad de Guadalajara, México», *Revista de Ciencias Ambientales*, **54**(2), págs. 1–19. doi: 10.15359/rca.54-2.1.
- [20] SIMA (2015), «Sistema Integral de Monitoreo Ambiental», <http://aire.nlgob.mx>.
- [21] TO, T., J. ZHU, D. STIEB, N. GRAY, I. FONG, L. PINAULT, M. JERRETT, A. ROBICHAUD, R. MÉNARD, A. VAN DONKELAAR *et al.* (2020), «Early life exposure to air pollution and incidence of childhood asthma, allergic rhinitis and eczema», *European Respiratory Journal*, **55**(2). doi: 10.1183/13993003.00913-2019.
- [22] ZHANG, Z., B. DONG, S. LI, G. CHEN, Z. YANG, Y. DONG, Z. WANG, J. MA y Y. GUO (2019), «Exposure to ambient particulate matter air pollution, blood pressure and hypertension in children and adolescents: a national cross-sectional study in China», *Environment international*, **128**, págs. 103–108. doi: 10.1016/j.envint.2019.04.036.

RESUMEN AUTOBIOGRÁFICO

Selene Berenice Prado Prado

Candidato para obtener el grado de
Ingeniería en Tecnología de Software

Universidad Autónoma de Nuevo León
Facultad de Ingeniería Mecánica y Eléctrica

Tesis:

MODELADO Y VISUALIZACIÓN DE RELACIONES ENTRE
CONTAMINANTES DEL AIRE Y SALUD PÚBLICA

Nací el 30 de Junio de 2000 en Monterrey, Nuevo León, soy la mayor de cuatro hijos. Mi familia está conformada por mi madre Lilia Prado López, mi padre Adan Alfaro Lerma, y mis hermanos: Angel Alejandro Prado Prado, Estrella Belen Prado Prado, y Genesis Adali Alfaro Prado.

Desde pequeña me han gustado las matemáticas, aprender cómo funcionan los sistemas computacionales, y leer.

Durante los primeros semestres de mi carrera descubrí la inteligencia computacional, un área que me encantó al instante, en especial su rama de ciencia de datos, rama en la que espero seguir desarrollándome.

Otra cosa que me apasiona es dibujar y pintar, actividades que estaban dentro de mí pero que se avivaron cuando inició la pandemia en el año 2020.